

NCBI Molecular Biology Resources

NCBI

The National Center for Biotechnology Information (NCBI)

- Created as a part of NLM in 1988
 - Establish public databases
 - Research in computational biology
 - Develop software tools for sequence analysis
 - Disseminate biomedical information
- Tools: BLAST(1990), Entrez (1992)
- GenBank (1992)
- Free MEDLINE (PubMed, 1997)
- Human genome (2001)

NCBI

GenBank: NCBI's Primary Sequence Database

RELEASE
81.0

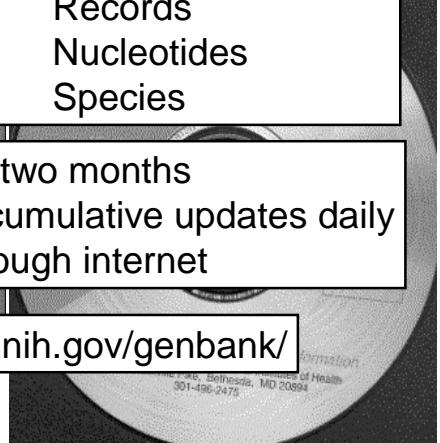
Release 129 April 2002	
16,769,983	Records
19,072,679,701	Nucleotides
110,000 +	Species

- full release every two months
- incremental and cumulative updates daily
- available only through internet

NCBI
National Library of Medicine, National Institutes of Health
8600 Rockville Pike, Bethesda, MD 20894

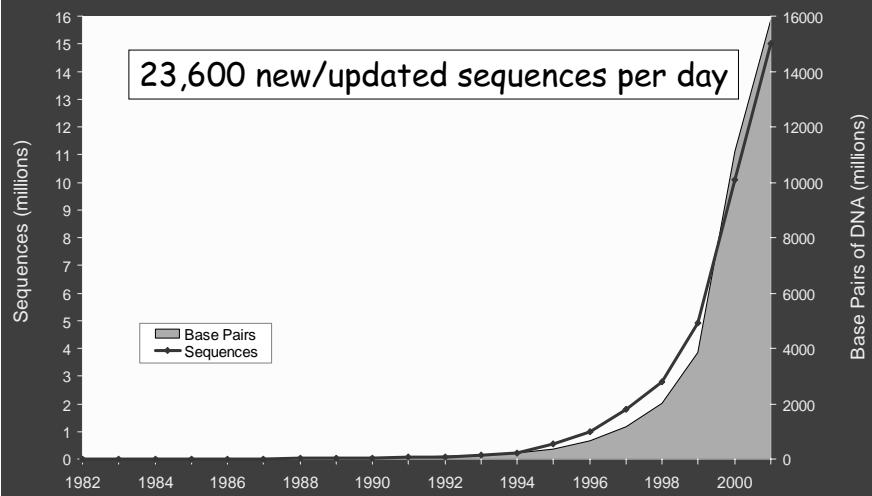
ftp://ftp.ncbi.nih.gov/genbank/

66.3 Gigabytes of data



NCBI

Growth of GenBank



RefSeq: NCBI's Derivative Sequence Database

- **Curated transcripts and proteins**
 - reviewed
 - human, mouse, rat, fruit fly, zebrafish, arabidopsis
- **Human model transcripts and proteins**
- **Assembled Genomic Regions (contigs)**
 - draft human genome
 - mouse genome
- **Chromosome records**
 - microbial
 - organelle

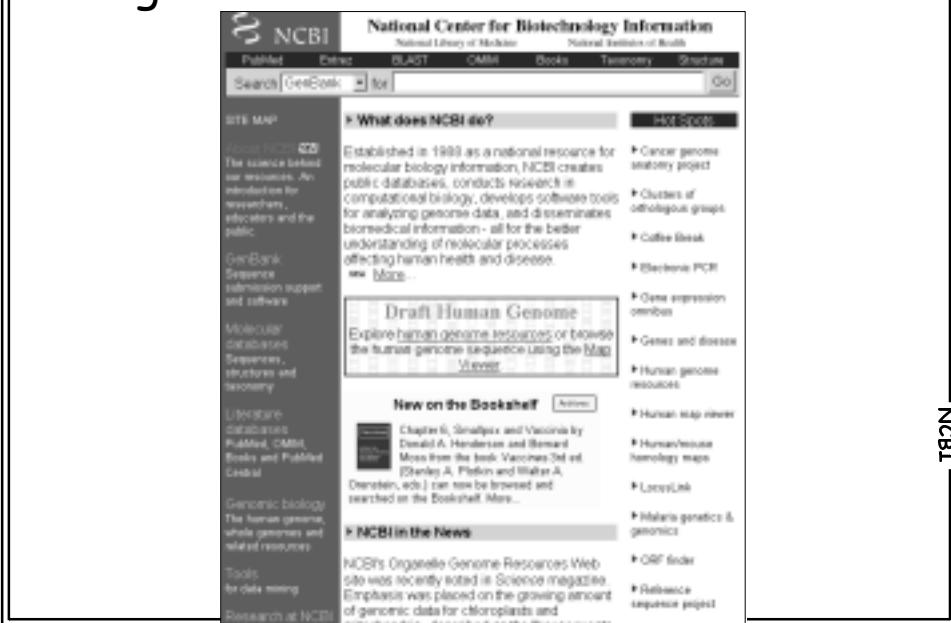
NCBI

RefSeq Resource

- **Genome Oriented Resource**
 - A sequence for each macromolecule Central Dogma: Chromosome, mRNA, preprotein, mature protein
 - Linked on a residue by residue basis
 - Objectively non-redundant and comprehensive
- **Curated Resource**
 - Authoritative source by genome
 - Derived from GenBank but corrected, merged, extended
 - Publicly distributed, Entrez Genomes Web site
- **Reagents for Genome Annotation and Analysis**
- **Substrate for Functional Genomics**

NCBI

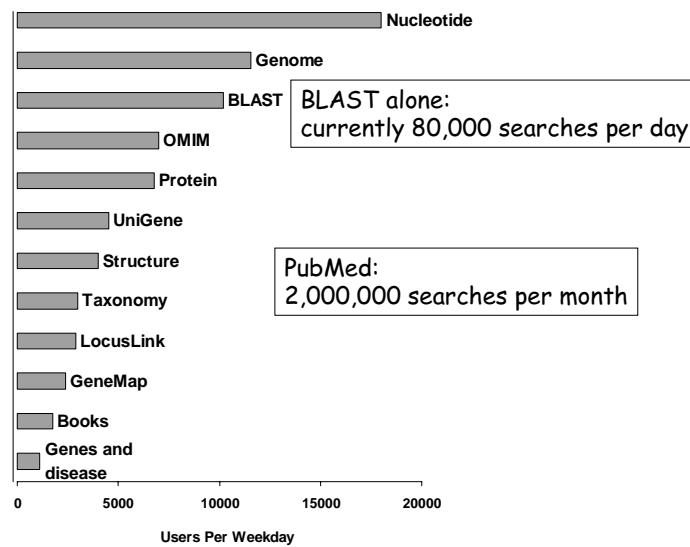
Integrated WWW Access: BLAST and Entrez



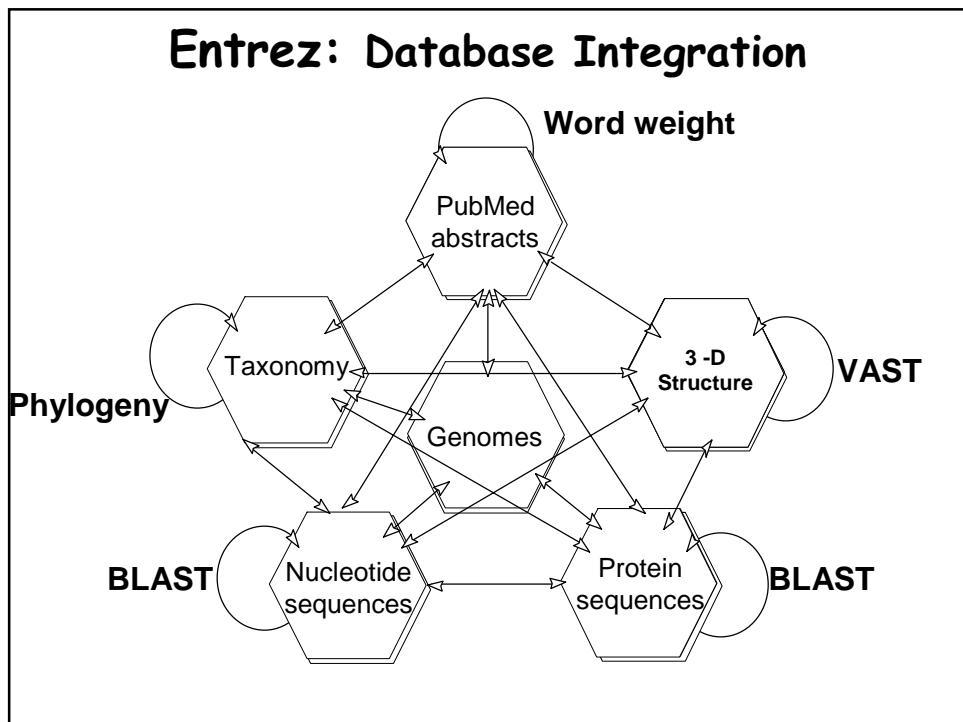
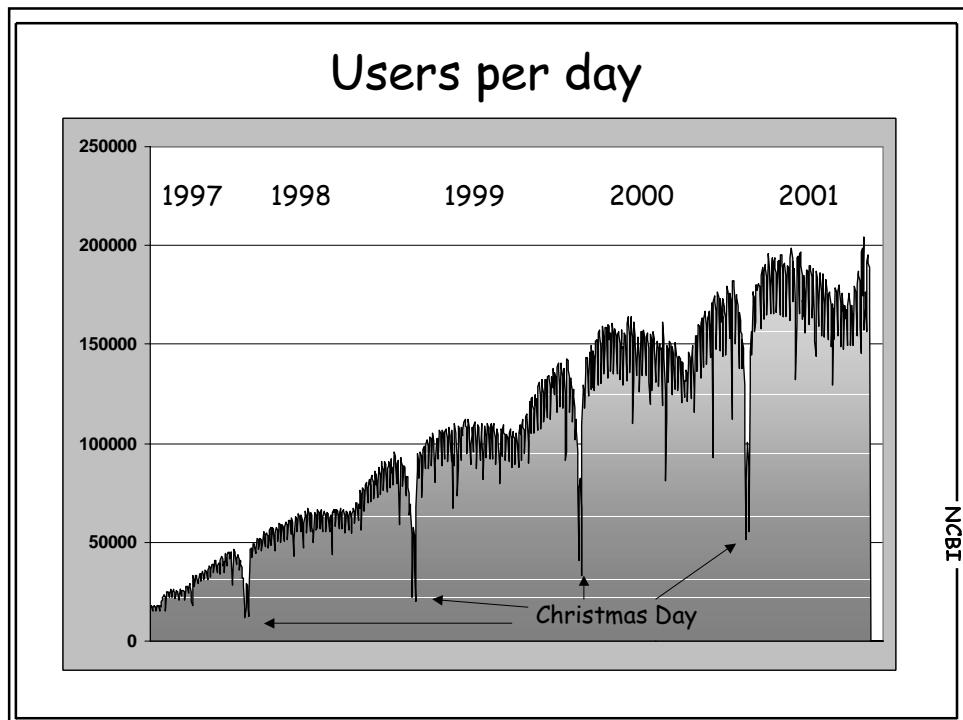
NCBI

Some Web Statistics

July 2001



NCBI



Basic Local Alignment Search Tool

NCBI

PubMed Entrez BLAST OMIM Taxonomy Structure

SITE MAP

BLAST info
BLAST overview
Frequently Asked Questions
New/Noteworthy
Receive e-mail with BLAST announcements
Description of BLAST Services
BLAST course
BLAST tutorial
BLAST references
URL API documentation
HTML format
PDF format
PostScript format

What's NEW in BLAST®

NEW March 5th 2002: New database linkouts from BLAST results. Results of a BLAST search will now link sequences from the BLAST results page to the NCBI LocusLink and UniGene databases. Links to additional databases coming soon

Nucleotide BLAST

- Standard nucleotide-nucleotide BLAST [blastn]
- MEGABLAST
- Search for short nearly exact matches

Protein BLAST

- Standard protein-protein BLAST [blastp]
- PSI- and PHI-BLAST
- Search for short nearly exact matches

Translated BLAST Searches

- Nucleotide query - Protein db [blastx]
- Protein query - Translated db [blastm]
- Nucleotide query - Translated db [tblastx]

Search for conserved domains

- Search the Conserved Domain Database using RPS-BLAST
- Search by domain architecture [DART]

Pairwise BLAST

- BLAST 2 Sequences

Genomic BLAST pages

- Human Genome
- Mouse Genome
- Rat Genome
- Fugu rubripes
- Zebrafish Genome
- Anopheles gambiae
- Arabidopsis thaliana
- Oryza sativa
- Other eukaryotes
- Microbial Genomes

Human Genome Resources Find a gene by ...

NCBI Home > Genomic Biology > Human

Search for

The Human Genome
A guide to online information resources

Web Resources

BLAST. Compare your sequence to the genome or its gene products.

Cytogenetics. A cytogenetic resource of FISH-mapped, sequence-tagged clones.

dbSNP. Database of SNPs and other genetic variations.

e-PCR. Check your sequence for STSs and view in genomic context.

GEO. Gene Expression Omnibus, a public repository for expression data.

HomoloGene. Putative homologies among human, mouse, rat, and zebrafish.

Homology Map. Blocks of conserved synteny between mouse and human.

LocusLink. Focal point for genes and associated information

Building an information infrastructure

A challenge facing researchers today is the ability to piece together and analyze the multitude of data currently being generated through the Human Genome Project. NCBI's Web site serves as an integrated, one-stop, genomic information infrastructure for biomedical researchers from around the world so that they may use this data in their research efforts. [More...](#)

Working Draft Analysis Published

- NLM Press Release
- NHGRI Press Release
- Interactive Tour of the Genome
- NCBI Genome Analysis Pipeline
- Nature (2/16/01) Human Genome Issue
- Science (2/16/01) Human Genome Issue

MapViewer tips and tricks

When browsing the genome using the new MapViewer, click on Display Settings to

OMIM. Guide to genes and inherited disorders maintained by JHU and collaborators.

RefSeq. Reference sequences of genomic contigs, mRNAs, and proteins.

SAGEmap. Gene expression results from SAGE tags mapped to sequences.

Sequencing. Summary of human genome sequencing progress.

MapViewer. Interactive viewer for genome maps, sequence, and genes.

UniGene. Organization of transcribed sequences stories into gene-based clusters.

UniSTS. A non-redundant collection of STSs with links to maps and sequence.

FTP Sites

dbSNP
LocusLink
protein Genomes
viral Genomes
RefSeq
UniGene

MapViewer

Below are three views of the BRCA2 locus using different display options. Click the image to see the full MapViewer display.

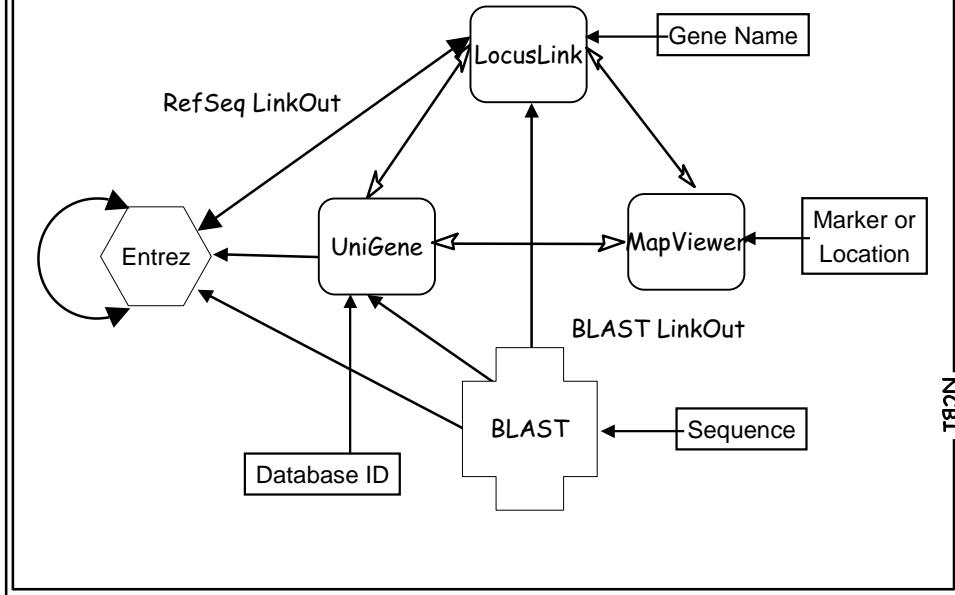
Other genomes

M. musculus
R. norvegicus
D. rerio
D. melanogaster
C. elegans
S. cerevisiae
organelles
bacteria
viruses

Frequently asked questions

What is a reference sequence?
How are gene symbols chosen for LocusLink and RefSeq?
How many inherited diseases have a known sequence?
How much of the human genome has been sequenced?
How were the genomic sequence contigs constructed?
How can I customize the MapViewer display?
What classes of genetic variation are included in dbSNP?
How can I deposit gene expression data in a public database?

Genome Resources Integration



Finding Human ESTs

Data		Score	E	Value
		(bits)		
Sequences producing significant alignments:				
Number of sequences	1	21	3e-06	U
	2	229	5e-58	U
Number of sequences	3	227	3e-56	U
	4	221	3e-56	U
Number of sequences	5	221	3e-56	U
	6	218	2e-55	U
Number of sequences	7	181	1e-52	U
	8	194	9e-52	U
Number of sequences	9	130	2e-50	U
	10	201	4e-50	U
Number of sequences	11	198	1e-49	U
	12	193	6e-48	U
Number of sequences	13	131	2e-47	U
	14	143	2e-45	U
Number of sequences	15	183	6e-45	U
	16	180	6e-44	U
Number of sequences	17	125	2e-43	U
	18	122	2e-43	U
Number of sequences	19	176	9e-43	U
	20	175	2e-42	U
Number of sequences	21	174	5e-42	U
	22	173	8e-42	U
Number of sequences	23	169	1e-40	U
	24	167	4e-40	U
Number of sequences	25	160	8e-40	U
	26	164	4e-39	U
Number of sequences	27	123	4e-39	U

MLH1 UniGene Cluster

UniGene Cluster Hs.57301 Homo sapiens

SEE ALSO

LocusLink: 4292
OMIM: 36
HomoloGene: 7301

SELECTED MC
organism, protein
H.sapiens: ORGANISM PROTEIN SIMILARITY
percent identity and length of a
P_000240.1 - MUTL PROTEIN HO

M.musculus: sp.Q5JK91 - MLH1 MOUSE DNA MISMATCH REPAIR PROTEIN MLH1 (MUTL PROTEIN HOMOLOGO 1)
R.norvegicus: sp.P97679 - MLH1_RAT DNA mismatch repair protein mlh1 (Mutl protein homolog 1)
A.thaliana: pir.T51620 - T51620 DNA mismatch repair protein MLH1 [imported] - Arabidopsis thaliana C
C.elegans: pir.T25389 - T25389 hypothetical protein Caenorhabditis elegans
E.coli: pir.PH0853 - PH0853 methyl-directed mismatch repair protein mutL - Escherichia coli
S.cerevisiae: sp.P38920 - MLH1 YEAST MUTL PROTEIN HOMOLOG 1 (DNA MISMATCH REPAIR PROTEIN MLH1)

MAPPING INFORMATION

Chromosome: 3
Genome View: Multiple mappings
OMIM Gene Map: Sp21.3
Whitehead map: WL-7345, Chr 3, YAC contig WC3
UniSTS entries: SHGC-12575 Genomic Context:

EXPRESSION INFORMATION

CDNA sources: 2 pooled primary tumors, one primary and brain; Stomach adenocarcinoma; Adenocarcinoma, cell line; melanoma, cell line; adenocarcinoma, cell line; melanoma, cell line; astrocytoma grade iv, cell line; lymphoma; blood; bone; brain; breast; breast lymphoma; carcinoma, cell line; choriocarcinoma; duodenal ..

mRNA/GENE SEQUENCES (6)

BC005866	Homo sapiens, Similar to mutL (E. coli) homolog 1 (colon cancer, nonpolyposis type 2), clone MGC:3965 IMAGE:2962831, mRNA, complete cds	[P]A
NM_000249	Homo sapiens mutL homolog 1, colon cancer, nonpolyposis type 2 (E. coli) (MLH1), mRNA	[P]A
U07418	Human DNA mismatch repair (hmlh1) mRNA, complete cds	[P]A
BC006850	Homo sapiens mutL (E. coli) homolog 1 (colon cancer, nonpolyposis type 2), clone MGC:5172 IMAGE:3451538, mRNA, complete cds	[P]A
U07343	Human DNA mismatch repair protein homolog (hMLH1) mRNA, complete cds	[P]A
BC005833	Homo sapiens, Similar to mutL (E. coli) homolog 1 (colon cancer, nonpolyposis type 2), clone MGC:2344 IMAGE:2962831, mRNA, complete cds	[P]A

EST SEQUENCES (10 of 278) [Show all ESTs]

BE884841	cDNA clone IMAGE:3911872 leiomyosarcoma	5' read [P]M
BG327257	cDNA clone IMAGE:4563987 renal cell adenocarcinoma	5' read [P]M
BF306562	cDNA clone IMAGE:4122836 rhabdomyosarcoma	5' read [P]M
B1256483	cDNA clone IMAGE:5122806 choriocarcinoma	5' read [P]M
BF795990	cDNA clone IMAGE:4343565 lymphoma, cell line	5' read [P]M
B1835081	cDNA clone IMAGE:5226825 pooled pancreas and spleen	5' read [P]M
BE539316	cDNA clone IMAGE:3451538 placenta	5' read [P]M
B1084102	cDNA clone IMAGE:5013904 epidermoid carcinoma, cell line	5' read [P]M
BG772547	cDNA clone IMAGE:4837570 testis, cell line	5' read [P]M
BG772733	cDNA clone IMAGE:4837559 testis, cell line	5' read [P]M

LocusLink MLH1

Click to Display mRNA-Genomic Alignments (spanning 57319 bps)

PUB	OMIM	UNIGENE	MAP	VAR	HOMOL	GDB	HGMD
ef	UCSC	PROTEOME	MCG				
Homo sapiens Official Gene Symbol and Name (HGNC)							
MLH1: mutL homolog 1, colon cancer, nonpolyposis type 2 (E. coli)							
LocusID: 4292							
Overview ?							
RefSeq Summary: MLH1 was identified as a locus frequently mutated in hereditary non-polyposis colon cancer (HNPCC). When cloned, it was discovered to be a human homolog of the E. coli mismatch repair gene mutS, consistent with the characteristic alterations in microsatellite sequences (RER+ phenotype) found in HNPCC.							
Protein Summary: MutL homolog 1, a mismatch repair protein							
Locus Type: gene with protein product, function known or inferred							
Product: mutL homolog 1							
Alternate Symbols: COCA2, HNPCC, hMLH1, HNPCC2							
Alias: mutL (E. coli) homolog 1 mutL (E. coli) homolog 1 (colon cancer, nonpolyposis type 2)							
Function: Submit GeneRIF (All Pubs)							
Phenotype: • Colorectal cancer, hereditary nonpolyposis, by inheritance • Leukemia • Muir-Torre family cancer syndrome • Turcot syndrome with glioblastoma							
GeneRIF: Gene References into Function: 11524701 • mutational analysis in HNPCC 11691925 • binds Bloom syndrome protein, nuclear localizat 11474654 • hereditary and somatic mutations in sporadic endometrial adenocarcinoma							
11809883 • hMutSalpha forms an ATP-dependent complex with hMutLalpha and hMutLbeta on DNA							
11748856 • Further characterization of the mutational spectrum of MLH1 gene in HNPCC families • PMG2-MLH1 protein interactions diminished by single nucleotide polymorphisms in HNPCC							

Links:

- pm** [PubMed](#)
- mv** [MapViewer](#)
- sv** [Sequence Viewer](#)
- ev** [Evidence Viewer](#)
- BL** [BLASTLink](#)

NCBI Reference Sequences (RefSeq)

Category: REVIEWED

mRNA: NM_000249
Protein: NP_000240 mutL homolog 1
Domains: Histidine kinase- score: 94
Histidine kinase-like ATPases score: 92
DNA mismatch repair protein. Also score: 513
known as the mutL/hexB/PMS1 family

GenBank: U07343
Source: ?

Category: NCBI Genome Annotation

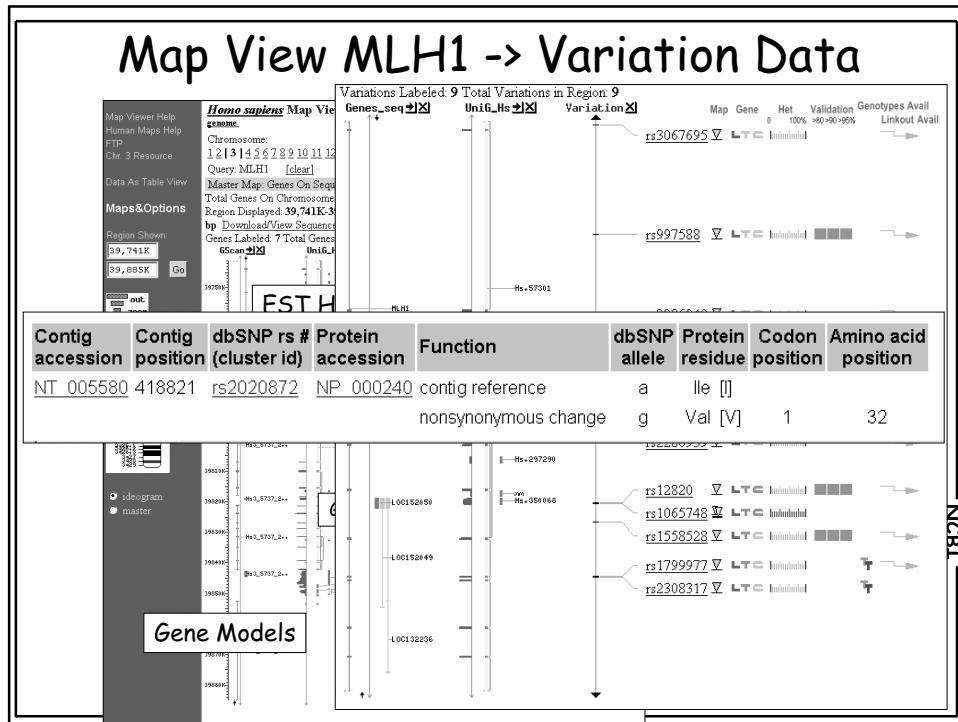
Genomic Contig: NT_005580

Annotated transcripts/proteins for this locus: sv mv ev

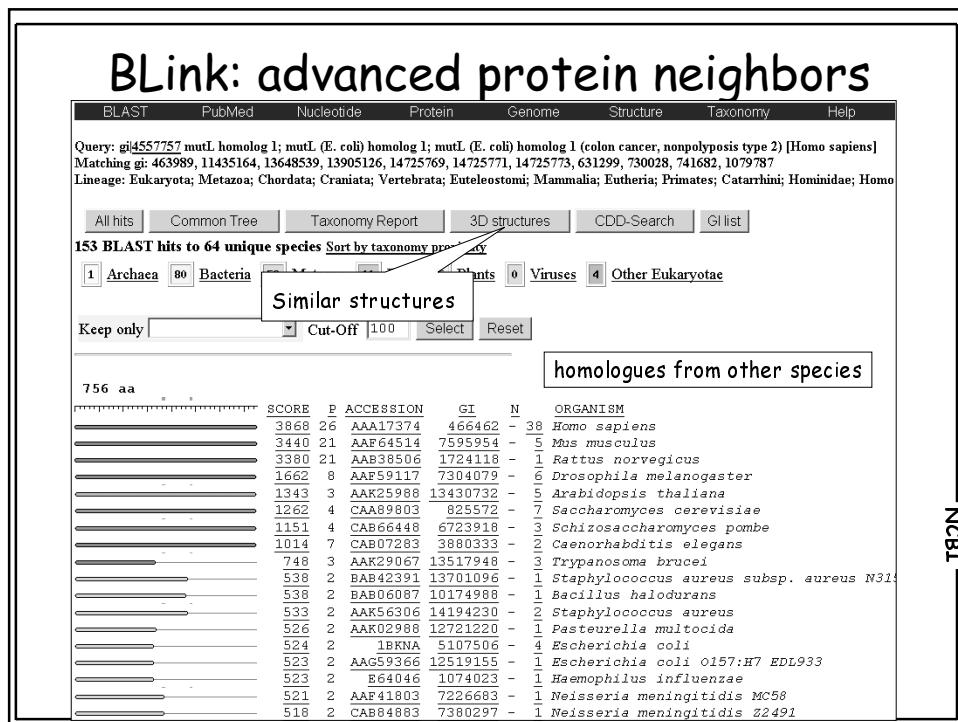
Evidence: supported by alignment with mRNA
NM_000249
NP_000240

BL

NCBI

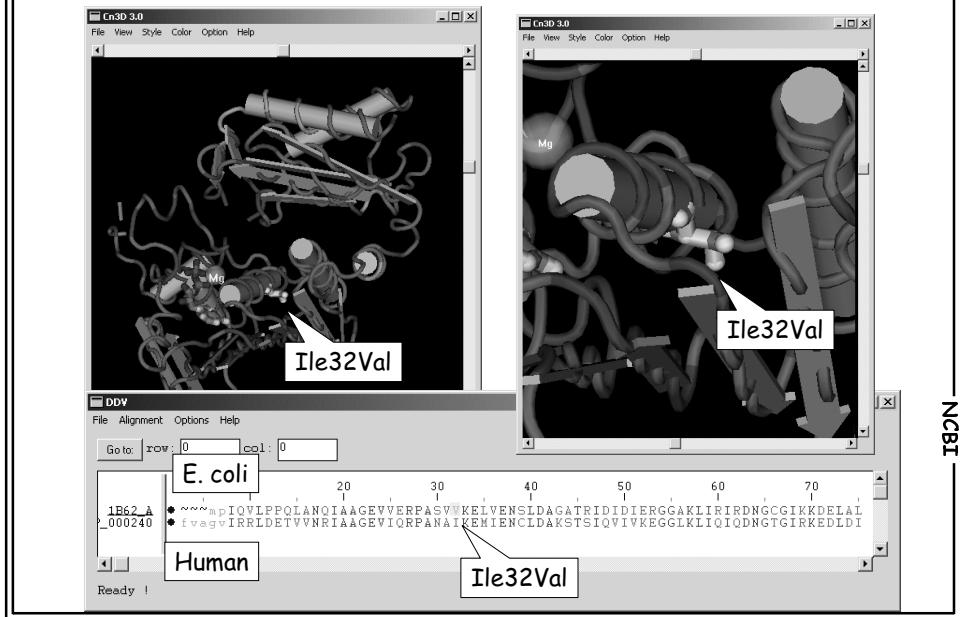


NCBI



NCBI

Cn3D: Finding a Modeling Template



NCB